**WHAT IS CLAIMED IS:**

1. A method for determining a language in which a document is created comprising the steps of:

    a) receiving at least one electronic document;

    b) identifying at least one character set encoding used in the at least one electronic document;

    c) determining whether the at least one character set encoding identifies a language in which the electronic document is created; and

    d) indicating the language in which the electronic document is created if a determination is made that the at least one character set encoding identifies the language in which the electronic document is created.

2. The method of claim 1, wherein the step of c) determining determines that the at least one character set encoding identifies at least two potential languages in which the electronic document is created.

3. The method of claim 2, further comprising the step of e) comparing at least one group of characters in the electronic document to predetermined groups of characters.

4. The method of claim 3, further comprising the step of f) detecting at least one identification for the at least one group of characters.

5. The method of claim 3, wherein the at least one group of characters is an n-gram.

6. The method of claim 4, wherein the at least one identification is a bit-flag.

7. The method of claim 4, further comprising the step of g) logically ANDing the at least one identification.

8. The method of claim 7, wherein the step of g) logically ANDing the at least one identification is repeated until a single identification is determined.

9. The method of claim 8, further comprising the step of h) indicating the language in which the electronic document is created.

10. The method of claim 9, further comprising the step of i) identifying a character set encoding for the language indicated.

11. A system for determining a language in which a document is created comprising:

    receiving means for receiving at least one electronic document;

    identifying means for identifying at least one character set encoding used in the at least one electronic document;

    determining means for determining whether the at least one character set encoding identifies a language in which the electronic document is created; and

    indicating means for indicating the language in which the electronic document is created if a determination is made that the at least one character set encoding identifies the language in which the electronic document is created.

12. The system of claim 11, wherein the determining means determines whether the at least one character set encoding identifies at least two potential languages in which the electronic document is created.

13. The system of claim 12, further comprising comparing means for comparing at least one group of characters in the electronic document to predetermined groups of characters.

14. The system of claim 13, further comprising detecting means for detecting at least one identification for the at least one group of characters.

15. The system of claim 13, wherein the at least one group of characters is an n-gram.

16. The system of claim 14, wherein the at least one identification is a bit-flag.

17. The system of claim 14, further comprising logical ANDing means for logically ANDing the at least one identification.

18. The system of claim 17, wherein the logically ANDing means logically ANDs the at least one identification until a single identification is determined.

19. The system of claim 18, further comprising language indicating means for indicating the language in which the electronic document is created.

20. The system of claim 19, further comprising character set encoding identifying means for identifying a character set encoding for the language indicated.

21. A system for determining a language in which a document is created comprising:

a receiving module that receives at least one electronic document;

an identifying module that identifies at least one character set encoding used in the at least one electronic document;

a determining module that determines whether the at least one character set encoding identifies a language in which the electronic document is created; and

an indicating module that indicates the language in which the electronic document is created if a determination is made that the at least one character set encoding identifies the language in which the electronic document is created.

22. The system of claim 21, wherein the determining module determines whether the at least one character set encoding identifies at least two potential languages in which the electronic document is created.

23. The system of claim 22, further comprising a comparing module that compares at least one group of characters in the electronic document to predetermined groups of characters.

24. The system of claim 23, further comprising a detecting module that detects at least one identification for the at least one group of characters.

25. The system of claim 23, wherein the at least one group of characters is an n-gram.

26. The system of claim 24, wherein the at least one identification is a bit-flag.

27. The system of claim 24, further comprising a logical ANDing module that logically ANDs the at least one identification.

28. The system of claim 27, wherein the logically ANDing module logically ANDs the at least one identification until a single identification is determined.

29. The system of claim 28, further comprising a language indicating module that indicates the language in which the electronic document is created.

30. The system of claim 29, further comprising a character set encoding identifying module that identifies a character set encoding for the language indicated.

31. A processor readable medium comprising processor readable code that causes a processor to determine a language in which a document is created, the processor readable medium comprising:

receiving code that causes a processor to receive at least one electronic document;

identifying code that causes a processor to identify at least one character set encoding used in the at least one electronic document;

determining code that causes a processor to determine whether the at least one character set encoding identifies a language in which the electronic document is created; and

indicating code that causes a processor to indicate the language in which the electronic document is created if a determination is made that the at least one character set encoding identifies the language in which the electronic document is created.

32. The medium of claim 31, wherein the determining code determines whether the at least one character set encoding identifies at least two potential languages in which the electronic document is created.

33. The medium of claim 32, further comprising comparing code that causes a processor to compare at least one group of characters in the electronic document to predetermined groups of characters.

34. The medium of claim 33, further comprising detecting code that causes a processor to detect at least one identification for the at least one group of characters.

35. The medium of claim 33, wherein the at least one group of characters is an n-gram.

36. The medium of claim 34, wherein the at least one identification is a bit-flag.

37. The medium of claim 34, further comprising logical ANDing code that causes a processor to logically AND the at least one identification.

5      38. The medium of claim 37, wherein the logically ANDing code logically ANDs the at least one identification until a single identification is determined.

39. The medium of claim 38, further comprising language indicating code that causes a processor to indicate the language in which the electronic document is created.

10      40. The medium of claim 39, further comprising character set encoding identifying code that causes a processor to identify a character set encoding for the language indicated.